# Leap to Petascale

## Developing Developer Tools
## TotalView and Blue Gene/Q

**May 23 2012**

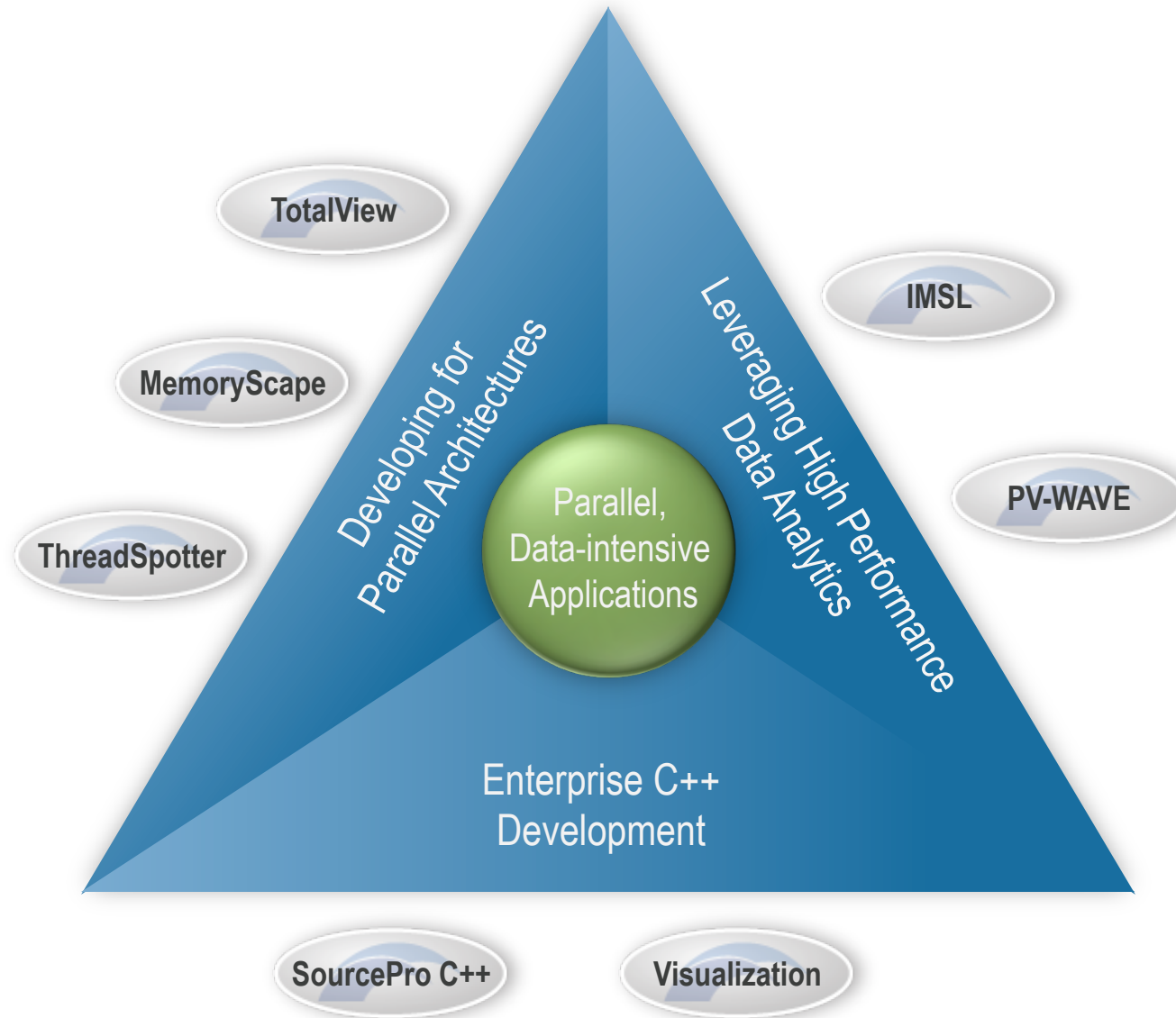Ed Hinkel,
Sr Sales Engineer
Rogue Wave Software

# Agenda

- **Who is Rogue Wave?**

- **Early Blue Gene Days with TotalView**

- **Blue Gene/Q Advancements**

- **Techniques for Debugging Challenges**

- **What's New with TotalView**
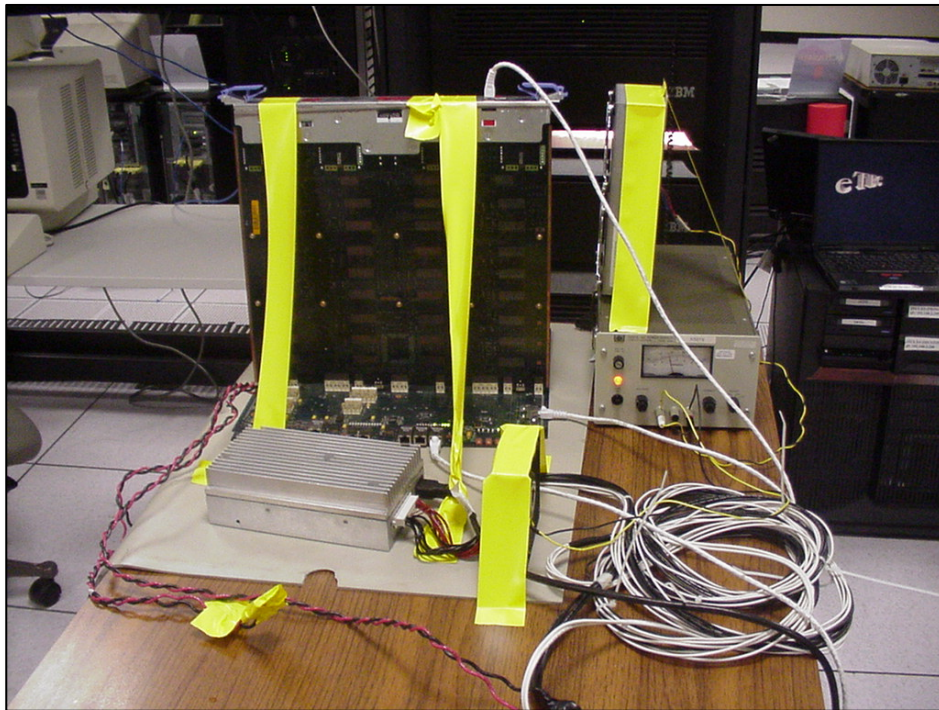
# Who is Rogue Wave Software?
# Solution Portfolio

# Early Blue Gene Days with TotalView

# TotalView Blue Gene Support

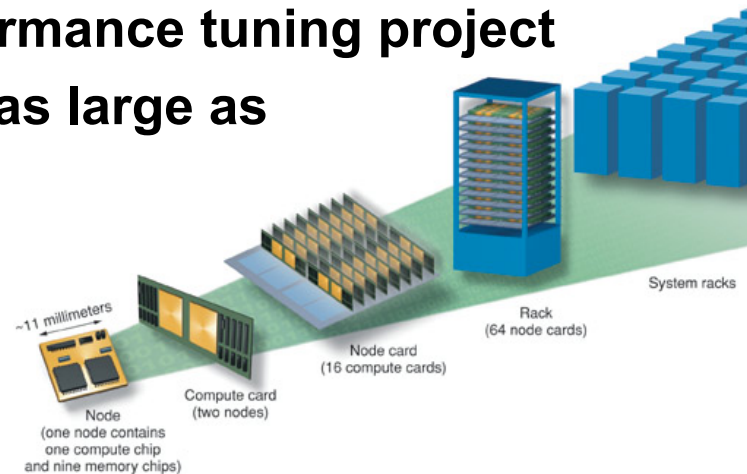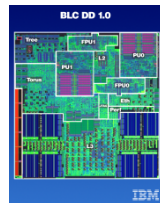- ## TotalView involvement started in 2003 on BG/L



IBM/TV BG/L development system

Gotta love that yellow duct tape!

ROGUE WAVE®
SOFTWARE

# TotalView Blue Gene/L Support

- **Support for Blue Gene/L since 2005**
- **Debugging interfaces developed via close collaboration with IBM (CIOD)**
- **Used on DOE/NNSA/LLNL's Blue Gene/L system containing 212 K cores**
  - **Heap memory debugging support added**
  - **Blue Gene/L scaling and performance tuning project**
  - **TotalView has debugged jobs as large as 32,768 processes**
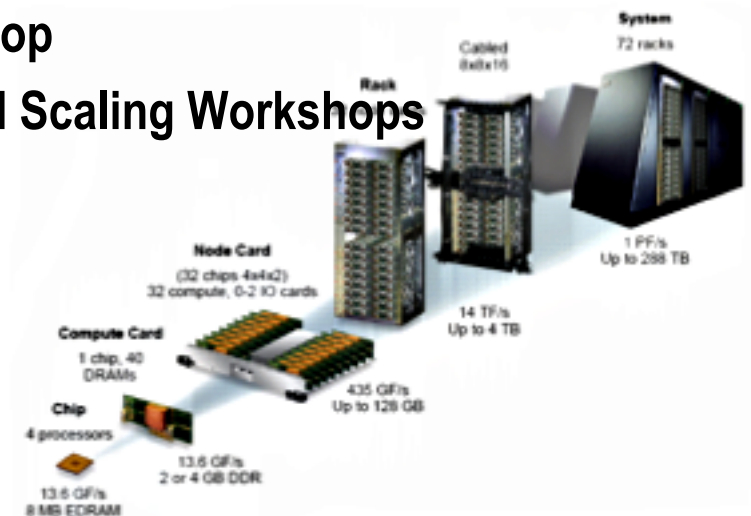


**Blue Gene/L work facilitated Blue Gene/P support**
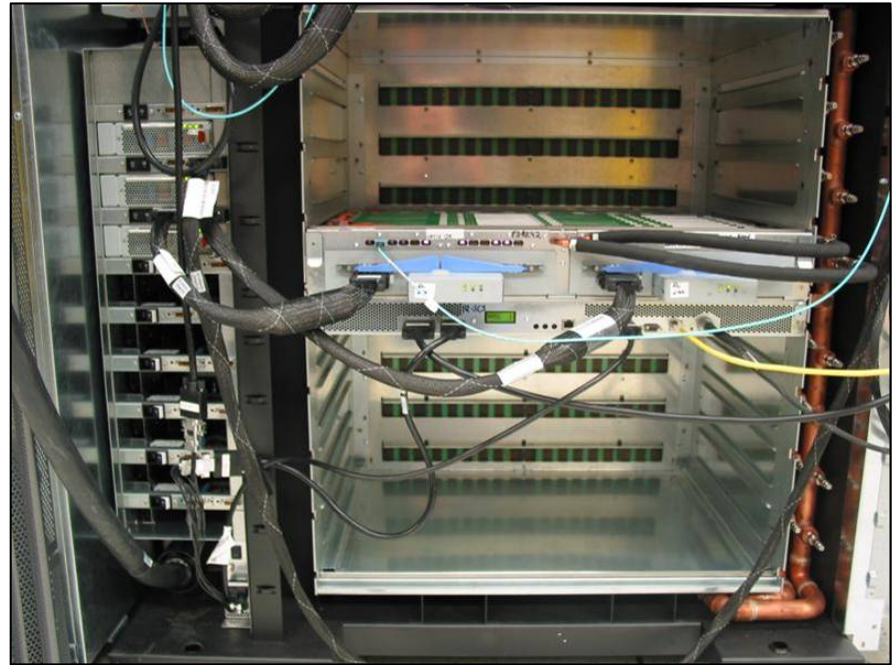
# TotalView Blue Gene/P Support

- **Continued close collaboration with IBM**

- **Currently running on several BG/P installations in Germany, France, the UK, and the US.**

- **Support for shared libraries, threads, and OpenMP**

- **TotalView has debugged jobs as large as 32,768**

- **Active workshop participation through the development**
  - **ANL's ALCF INCITE Performance Workshop**
  - **Jülich's Blue Gene/P Porting, Tuning, and Scaling Workshops**

# TotalView Blue Gene/Q Support

- **Porting TotalView began in June 2011**

- **Access to Q32 at IBM began in August**

- **Basic debugging operations in October**

- **Used in Synthetic Workload Testing in December**
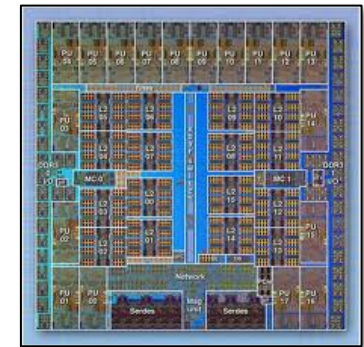
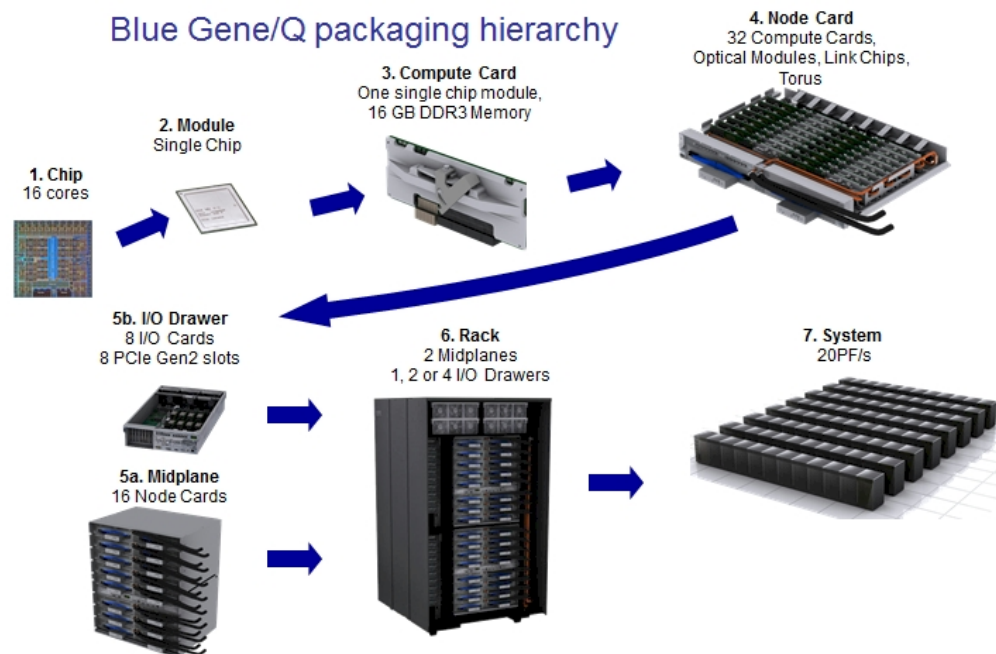- **Fully functional in March 2012**



**IBM's Q32**

# TotalView Blue Gene/Q Support (cont)

- **Thanks to the ongoing collaboration with IBM and the BG Kernel Team, early access versions of TotalView for BG/Q isavailable**

- **Argonne National Laboratory**

- **Lawrence Livermore National Laboratory**



## Blue Gene/Q packaging hierarchy

1. Chip
16 cores

2. Module
Single Chip

3. Compute Card
One single chip module,
16 GB DDR3 Memory

4. Node Card
32 Compute Cards,
Optical Modules, Link Chips,
Torus

5b. I/O Drawer
8 I/O Cards
8 PCIe Gen2 slots

5a. Midplane
16 Node Cards

6. Rack
2 Midplanes
1, 2 or 4 I/O Drawers

7. System
20PF/s

# Blue Gene/Q Advancements

# with

# TotalView

# TotalView on BG/Q Support

- **BG/Q TotalView is as functional as BG/P TotalView**
  - MPI, OpenMP, pthreads, hybrid MPI+threads
  - C, C++, Fortran, assembler; IBM and GNU compilers
  - Basics: source code, variables, breakpoints, watchpoints, stacks, single stepping, read/write memory/registers, conditional breakpoints, etc.
  - Memory debugging, message queues, binary core files, etc.

- **PLUS, features unique to BG/Q TotalView**
  - QPX (floating point) instruction set and register model
  - Fast compiled conditional breakpoints and watchpoints
  - Asynchronous thread control

- **Working with IBM on debugging interfaces for TM/SE regions**
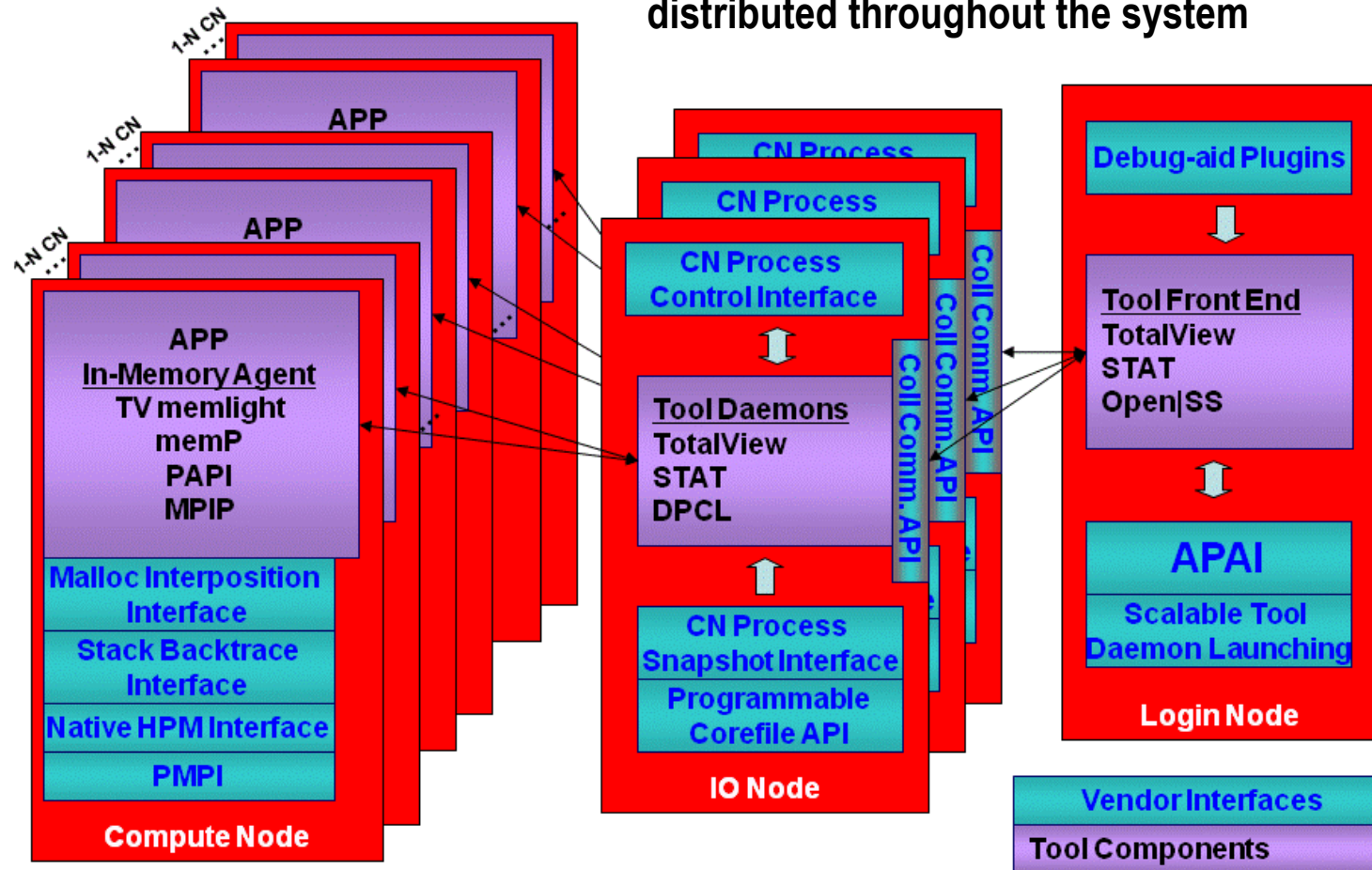  - TM == transactional memory; SE == speculative execution

**ROGUE WAVE®**
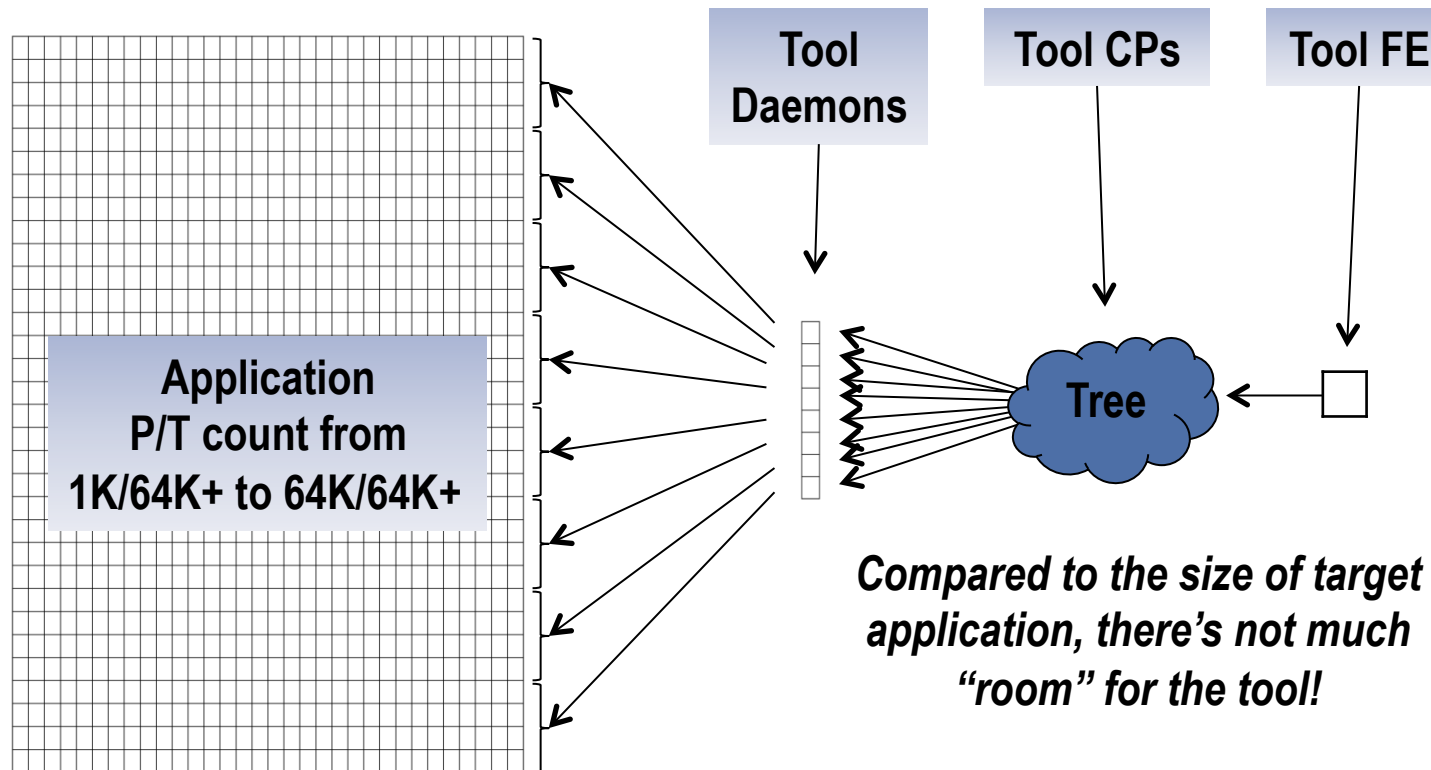SOFTWARE

# Advanced TotalView Features on BG/Q

- **Asynchronous thread control**
  - A feature on Linux, and other TotalView platforms, ported to BG/Q
  - Allows you to individually control the execution of threads
  - Run and halt individual threads
  - Single-step a group of threads in lockstep
  - Hold and release the execution of individual threads
  - Create stop-thread and thread barrier breakpoints
- **Fast compiled conditional breakpoints and watchpoints**
  - A feature on AIX and other TotalView platforms, ported to BG/Q
  - Conditional breakpoints and watchpoints execute in as little as 7 μsecs
  - Conditional expressions are compiled and dynamically patched into the process
  - Evaluation is performed by the triggering thread, in parallel

# Blue Gene Code Development Tools Interface (CDTI)

**Hierarchical infrastructure components are distributed throughout the system**

# Tool Challenges

**Application P/T count from 1K/64K+ to 64K/64K+**

**Tool Daemons**

**Tool CPs**

**Tool FE**

**Tree**

*Compared to the size of target application, there's not much "room" for the tool!*

**One rack of BG/Q: 1K CNs, 16K Cores, 64K HW Threads**

**A "generous" 128:1 CN:ION ratio: 8 IONs**

**A "beefy" FEN P7, 3 GHz+, 32 GB+**

# Overcoming High CN:ION Ratios

- **On BG/Q, at a given ratio, on *each* IO node, tool daemons may be responsible for up to**

| CN:ION | Processes | Threads |
|--------|-----------|---------|
| 64:1 | 4,096 | 20,480 |
| 128:1 | 8,192 | 40,960 |
| 256:1 | 16,384 | 81,920 |
| 512:1 | 32,768 | 163,840 |

- **But each IO node has**
  - **1.6 GHz A2 17 core processor (not too swift)**
  - **16 GB (limited physical memory)**

# Where to put the "weight" of the debugger?

- **Most of the "weight" of the debugger is in the symbol table**

- **Real-world applications are huge and complex**

- **A recently analyzed mission critical application revealed**
  - **1.5 million function definitions**
  - **16 million line number definitions**
  - **DWARF symbol information in excess of 2 GB**
  - **100s or 1000 of shared libraries**

- **You don't want to be big in the back end!**

- **And nothing too compute intensive either**

# TotalView's Architecture

- **Extremely lightweight back-end daemon processes**
  - Small footprint plus a few hundred bytes per CN process or thread
  - Each daemon can handle thousands of processes and threads
  - The daemons do not store the symbol table!
- **The "weight" of the debugger is on the front-end node**
- **Symbol tables are indexed and stored on the FEN**
  - Debugger has exactly *one* copy of the symbol table for each image file
  - Symbol tables are *shared* across all processes and thread
  - Aggregate memory consumption is minimal

- **Well suited to Blue Gene!**

# There's still the high P/T count per IO node problem

- **Process and threads counts per IO node are still high!**

- **What to do about that?**

- **"Divide and conquer"**
  - Place a small number of daemons on the ION
  - We do have 17 cores we can use

- **Unlike CIOD on BG/L&P, CDTI on BG/Q can operate in parallel**
  - There's one CDTI debug channel per compute node

# Solution: TotalView/MRNet Trees on the IO Nodes



| 0-7 |
| 8-15 |
| 16-23 |
| 24-31 |
| 32-39 |
| 40-47 |
| 48-55 |
| 56-63 |
| 64-71 |
| 72-79 |
| 80-87 |
| 88-95 |
| 96-103 |
| 104-111 |
| 112-119 |
| 120-127 |

1 CDTI channel per CN

| tvdsvr0 |
| tvdsrv1 |
| tvdsvr2 |
| tvdsvr3 |
| tvdsvr4 |
| tvdsvr5 |
| tvdsvr6 |
| tvdsvr7 |
| tvdsvr |
| tvdsvr9 |
| tvdsvr10 |
| tvdsvr11 |
| tvdsvr12 |
| tvdsvr13 |
| tvdsvr14 |
| tvdsvr15 |

MRNet CP

MRNet Tree

*Instead of one daemon managing all 128 CNs*

*The MRNet Commnode Process connects the daemons to the rest of the tree*

**128 CNs**

**1 ION**

ROGUE WAVE
SOFTWARE

# TotalView on BG/Q Support

- **BG/Q TotalView is as functional as BG/P TotalView**
  - MPI, OpenMP, pthreads, hybrid MPI+threads
  - C, C++, Fortran, assembler; IBM and GNU compilers
  - Basics: source code, variables, breakpoints, watchpoints, stacks, single stepping, read/write memory/registers, conditional breakpoints, static/dynamic executables, etc.
  - Memory debugging, message queues, binary core files, etc.
- **PLUS, advanced BG/Q TotalView features**
  - QPX (floating point) instruction set and register model
  - Fast compiled conditional breakpoints and watchpoints
  - Asynchronous thread control
- **Working with IBM on debug interfaces for TM/SE regions**
  - TM == transactional memory; SE == speculative execution

# Advanced BG/Q TotalView Features

- ## Asynchronous thread control
  - A TotalView feature on Linux and other platforms, ported to BG/Q
  - Allows you to individually control the execution of threads
  - Run and halt individual threads
  - Single-step a group of threads in lockstep
  - Hold and release the execution of individual threads
  - Create stop-thread and thread barrier breakpoints
- ## Fast compiled conditional breakpoints and watchpoints
  - A TotalView feature on AIX and other platforms, ported to BG/Q
  - Conditional breakpoints and watchpoints execute
    in as little as 7 $\mu$secs
  - Conditional expressions are compiled and dynamically patched into the process
  - Evaluation is performed in parallel by the triggering thread

ROGUE WAVE®
SOFTWARE

# TotalView Availability

- ## TotalView on Blue Gene/Q Today
    - LLNL has it up and running on rzuseq, and is using it to debug applications.
    - IBM is using it internally for debugging and testing.
    - It's installed on IBM's Blue Gene On Demand Center Q32 (if anyone has access to that system).

- ## TotalView On Blue Gene at Argonne
    - 1024 Tokens (BG/P)
    - Research license is available with 65,536 tokens

ROGUE WAVE®
SOFTWARE

# TotalView on VEAS!

# TotalView on VEAS!

# Techniques for Debugging Challenges

ROGUE WAVE®
SOFTWARE
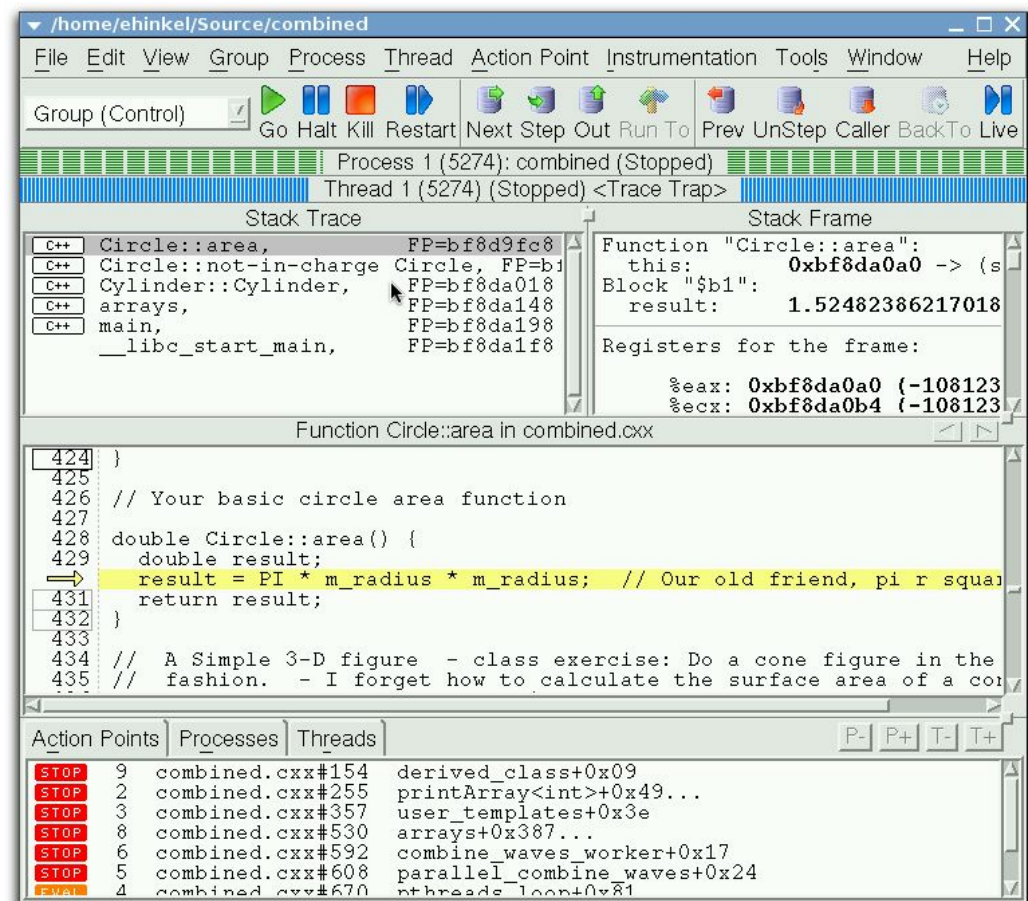
# What is TotalView?

## A comprehensive debugging solution for demanding parallel and multi-core applications

- **Wide compiler & platform support**
  - C, C++, Fortran 77 & 90, UPC
  - Unix, Linux, OS X

- **Handles Concurrency**
  - Multi-threaded Debugging
  - Multi-process Debugging

- **Integrated Memory Debugging**

- **Reverse Debugging available**

- **Supports Multiple Usage Models**
  - Powerful and Easy GUI – Highly Graphical
  - CLI for Scripting
  - Long Distance Remote Debugging
  - Unattended Batch Debugging



**ROGUE WAVE**
**SOFTWARE**

# Debugging Complex Codes

- **Mechanize**

- **Minimize**

- **Visualize**

- **… and Don't Forget the Memory**

# Mechanize

## Extended Automation Capabilities



|

# Automated Debugging

## TVscript

- **Non-Interactive Batch Debugging –**
  - **Work in the "main" batch queue**
  - **Don't have to baby-sit job waiting on it to run**
  - **Use scripting to perform checks that would be tedious to do by hand**
  - **Verification through automated processes (nightly build and test)**

## TTF and C++View

- **Automatic Transformation of Data –**
  - **Simplify interactive (and scripted) debugging**
  - **Perform validation/sanity checking of large datasets**
  - **Comparative debugging**
  - **Allows you to focus on troubleshooting your program**

ROGUE WAVE®
SOFTWARE

# Non-Interactive Batch Debugging with TVscript

- **Run multiple debugging sessions without the need for recompiling, unlike with printf**
- **TVscript syntax:**

  **tvscript [ options ] [ filename ] [ -a program_args ]**
- **More complex actions-to-events are possible, utilizing TCL within a CLI file**
- **TVscript lets you define what events to act on, and what actions to take**

<table>
<tr><td colspan="2"><strong>TVscript uses a simple, Event/Action interface</strong></td></tr>
<tr><td><strong>Typical Events</strong></td><td><strong>Typical Actions</strong></td></tr>
<tr><td>

- Action_point
- Any_memory_event
- Guard_corruption error

</td><td>

- Display_backtrace [-level *level-num*]
- List_leaks
- Save_memory
- Print [-slice {*slice_exp*] {*variable* | *exp*}

</td></tr>
</table>

29

ROGUE WAVE®
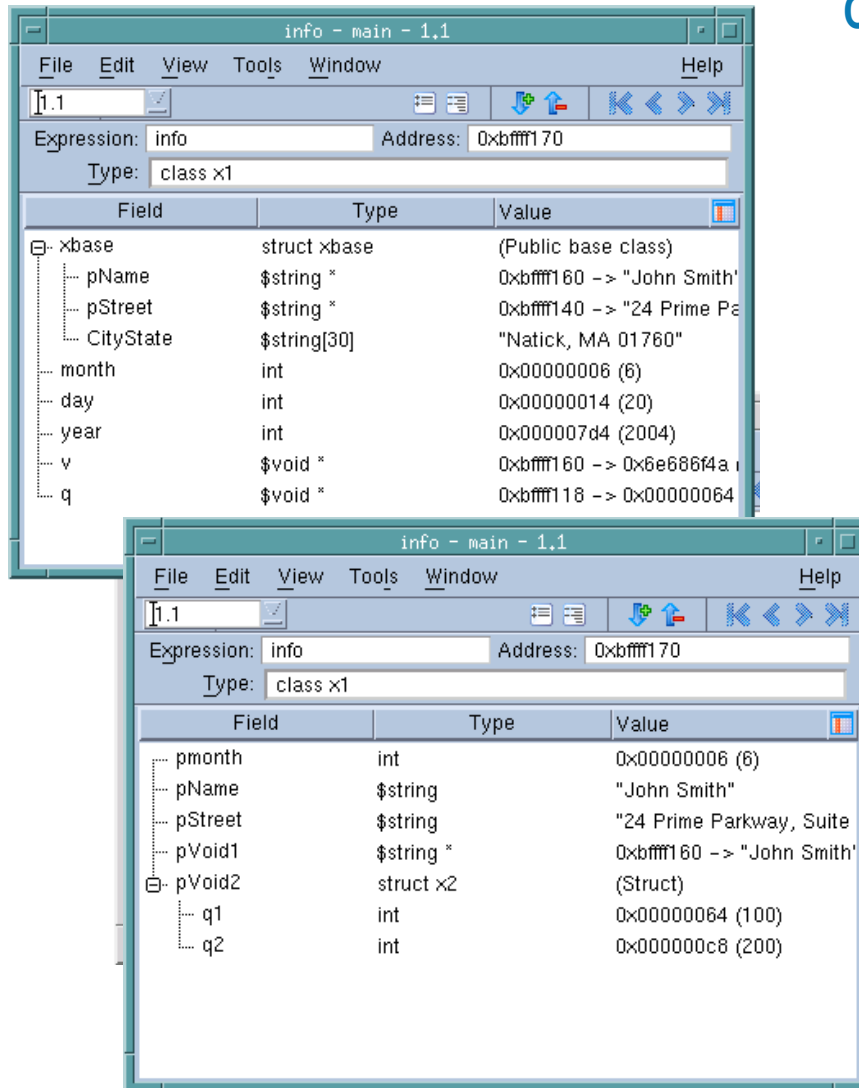SOFTWARE

# Unattended Debugging with Tvscript

**Example**

The following tells tvscript to report the contents of the *foreign_addr*
structure each time the program gets to line 85
**-create_actionpoint "#85=>print foreign_addr"**

Typical output sample with tvscript:

```
!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!
! Print
!
! Process:
!    ./server (Debugger Process ID:  1, System ID:  12110)
! Thread:
!    Debugger ID:  1.1, System ID:  3083946656
! Time Stamp:
!    06-26-2008 14:04:09
! Triggered from event:
!    actionpoint
! Results:
!    foreign_addr = {
!       sin_family = 0x0002 (2)
!       sin_port = 0x1fb6 (8118)
!       sin_addr = {
!          s_addr = 0x6658a8c0 (1717086400)
!       }
!       sin_zero = ""
!    }
!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!
```

ROGUE WAVE®
S O F T W A R E

# Creating Type Transformations



**Customize your own Transformations**

In $HOME/.tvdrc:

::TV::TTF::RTF::build_struct_transform {

    name    {^class x1$}

    members {

        { pmonth    { month } }

        { pName    { xbase upcast { * pName    } } }

        { pStreet   { xbase upcast { * pStreet } } }

        { pVoid1   { "$string *"   cast v     } }

        { pVoid2   { * { "class x2 *"   cast q } } }

    }

}

**Meta Language:**
**{member}**
**{* expr}**
**{expr . Expr}**
**{expr -> expr}**
**{datatype case expr}**
**{baseclass upcast expr}**

ROGUE WAVE®
S O F T W A R E

# C++View

- **C++View is an easy way to customize TotalView's display of object data.**

- **How does it work?**
  - User writes short display functions within their program
  - TotalView uses these functions to simplify the display of data when the user explores their data within that program
  - C++View transforms are easy to define
  - Great for collaborative codes (transforms can be distributed with the program)
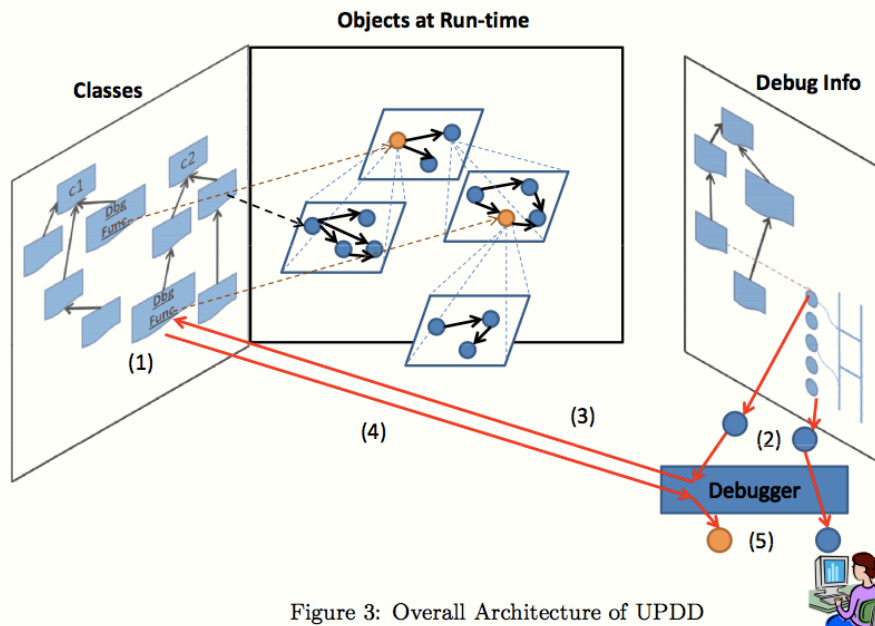
- **Benefit: Easier for scientists and developers to work with complex applications**

Figure 3: Overall Architecture of UPDD

Developers can now write display and analysis functions for their C++ classes that are invoked whenever an object is inspected interactively in the debugger.

# C++View

- **C++View is a simple way for you to define type transformations**
  - **Simplify complex data**
  - **Aggregate and summarize**
  - **Check validity**
- **Transforms**
  - **Type-based**
  - **Compose-able**
  - **Automatically visible**
- **Code**
  - **C++**
  - **Easy to write**
  - **Resides in target**
  - **Only called by TotalView**
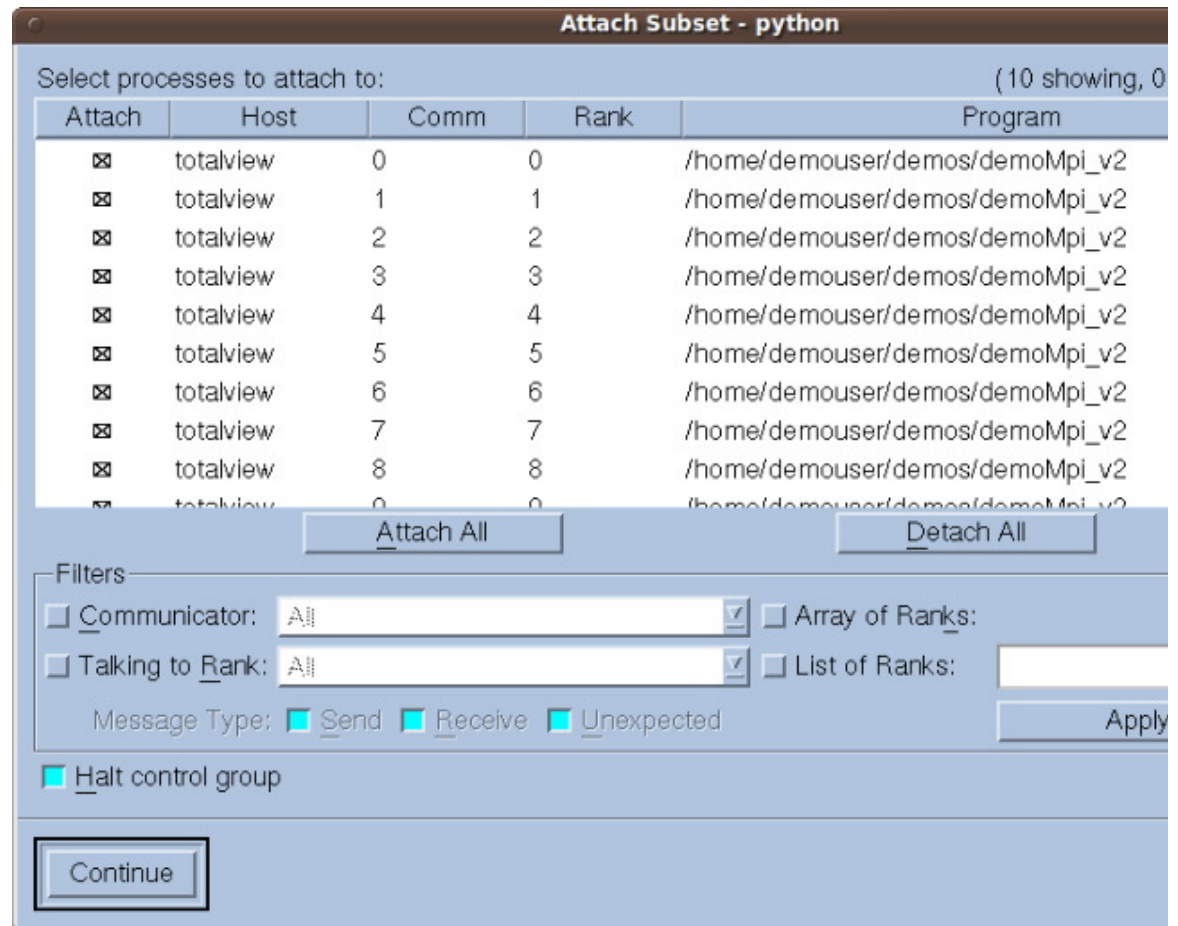
# Minimize

## Reduce the Scope of Effort

# Subset Debugging
## With TotalView
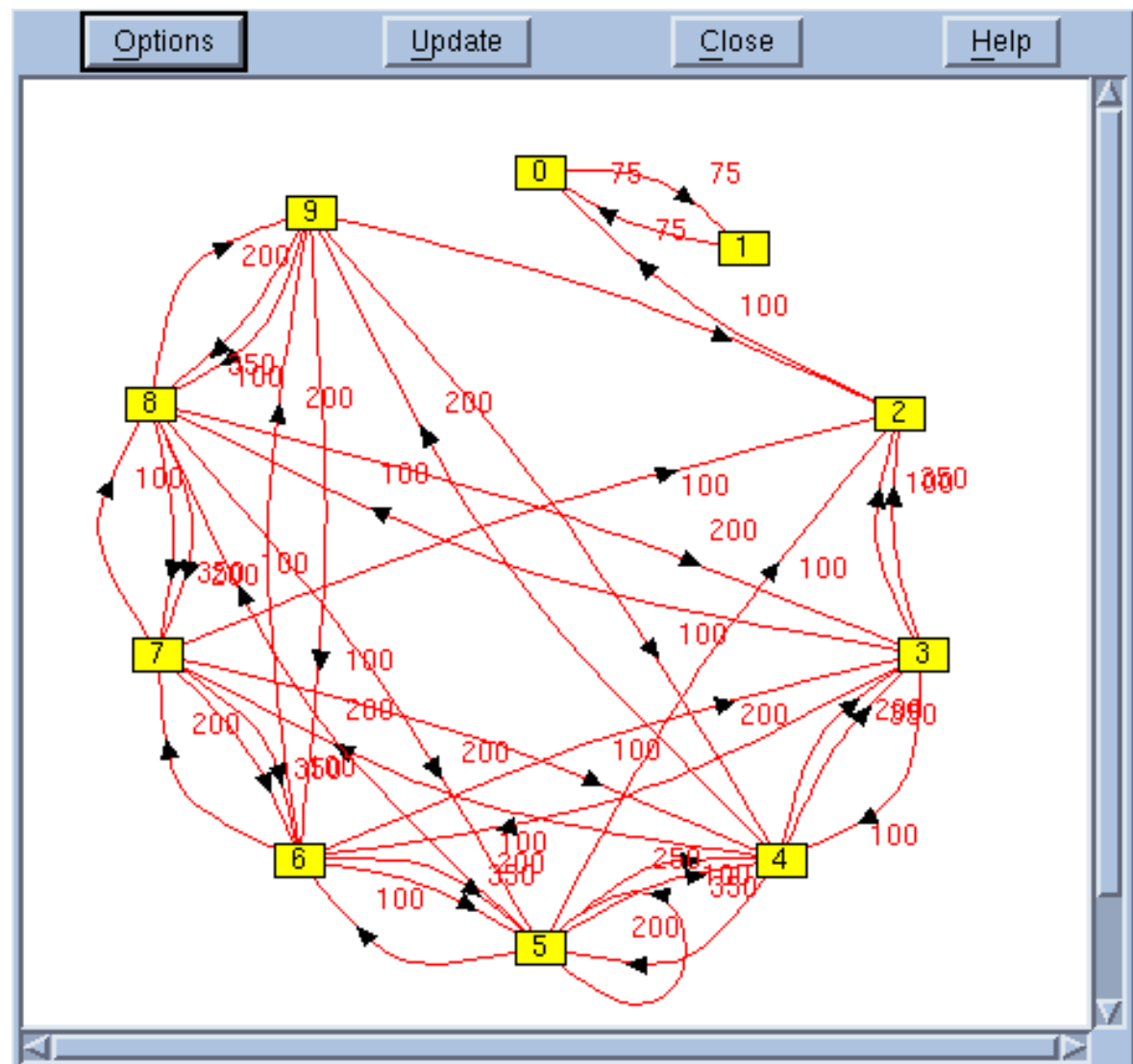
# Subset Attach

## You need not be attached to the entire job



- **You can be attached to different subsets at different times through the run**

- **You can attach to a subset, run till you see trouble and then 'fan out' to look at more processes if necessary.**

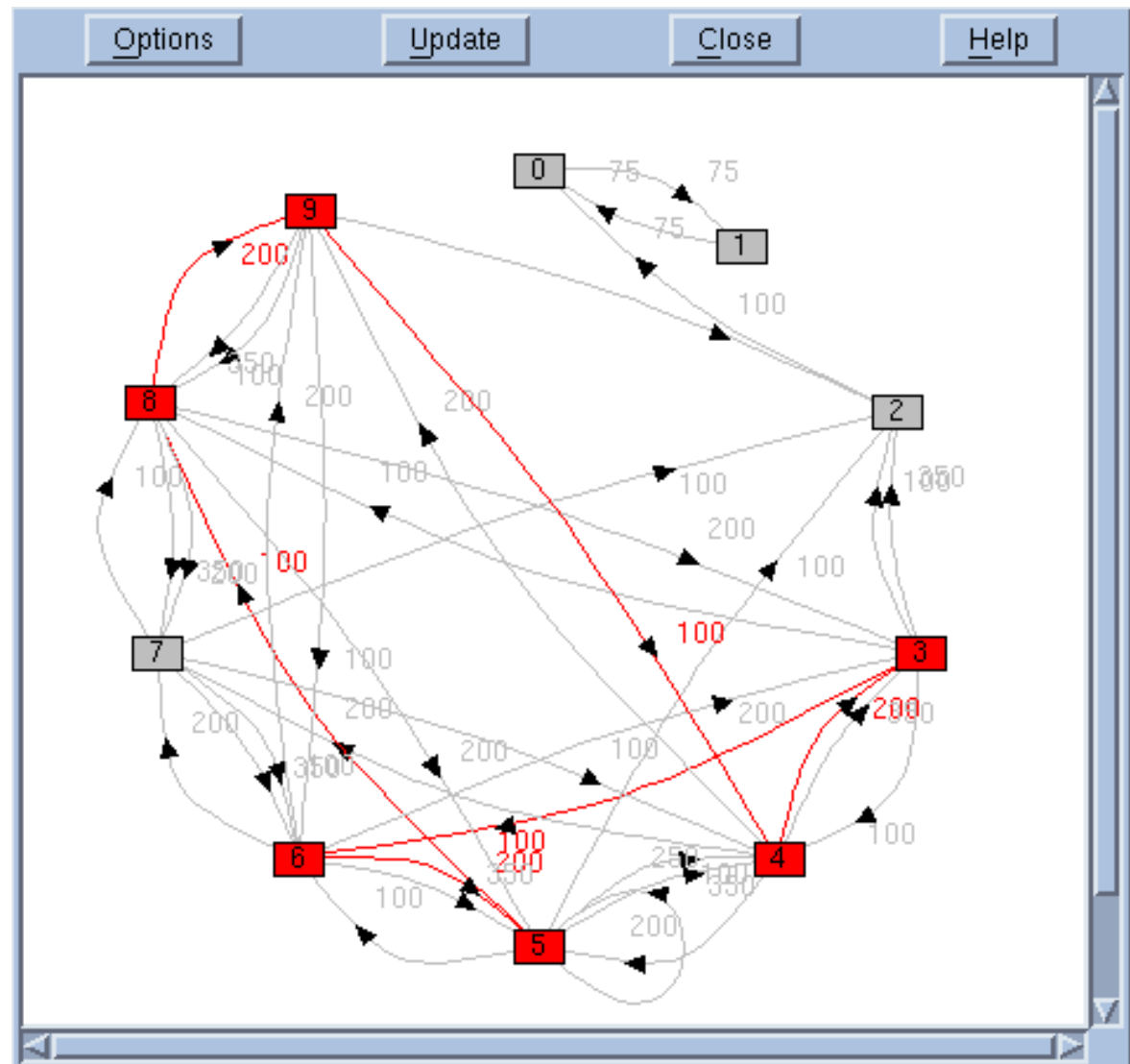- **This greatly reduces overhead**

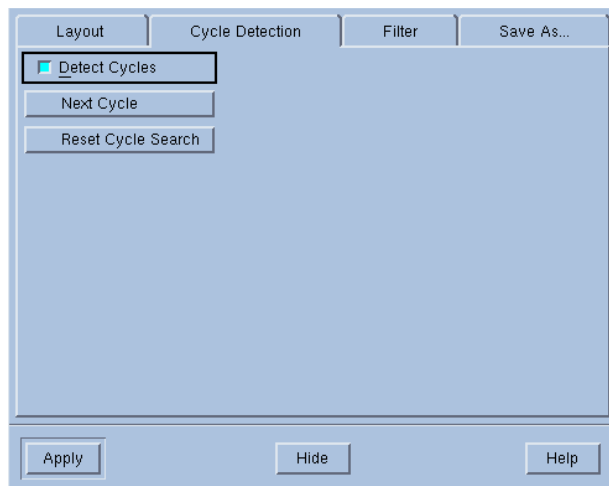- **It also reduces license size requirements**

# Message Queue Graph

- **Hangs & Deadlocks**
- **Pending Messages**
  - **Receives**
  - **Sends**
  - **Unexpected**
- **Inspect**
  - **Individual entries**
- **Patterns**

# Message Queue Debugging

- **Filtering**
  - **Tags**
  - **MPI Communicators**
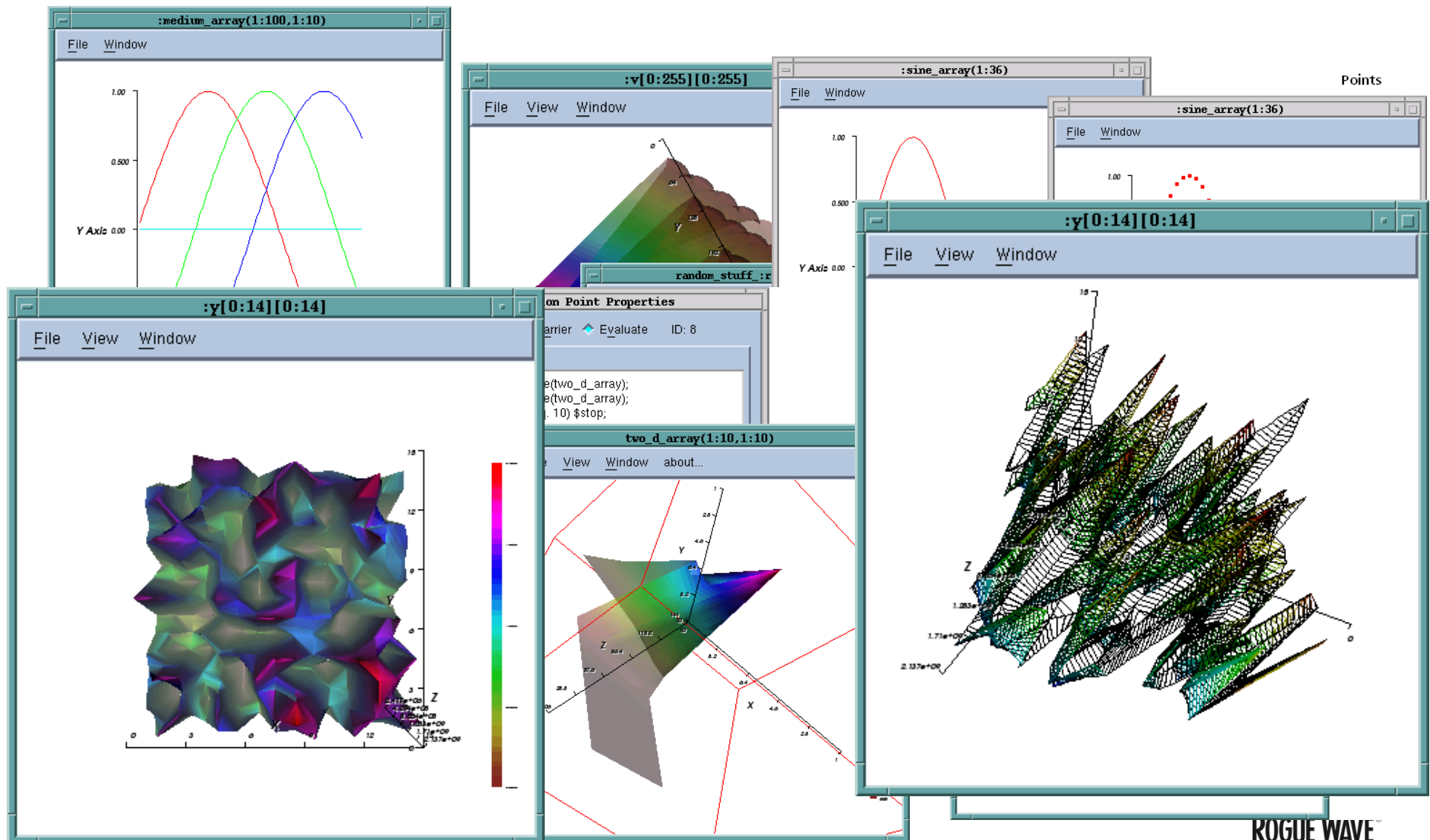- **Cycle detection**
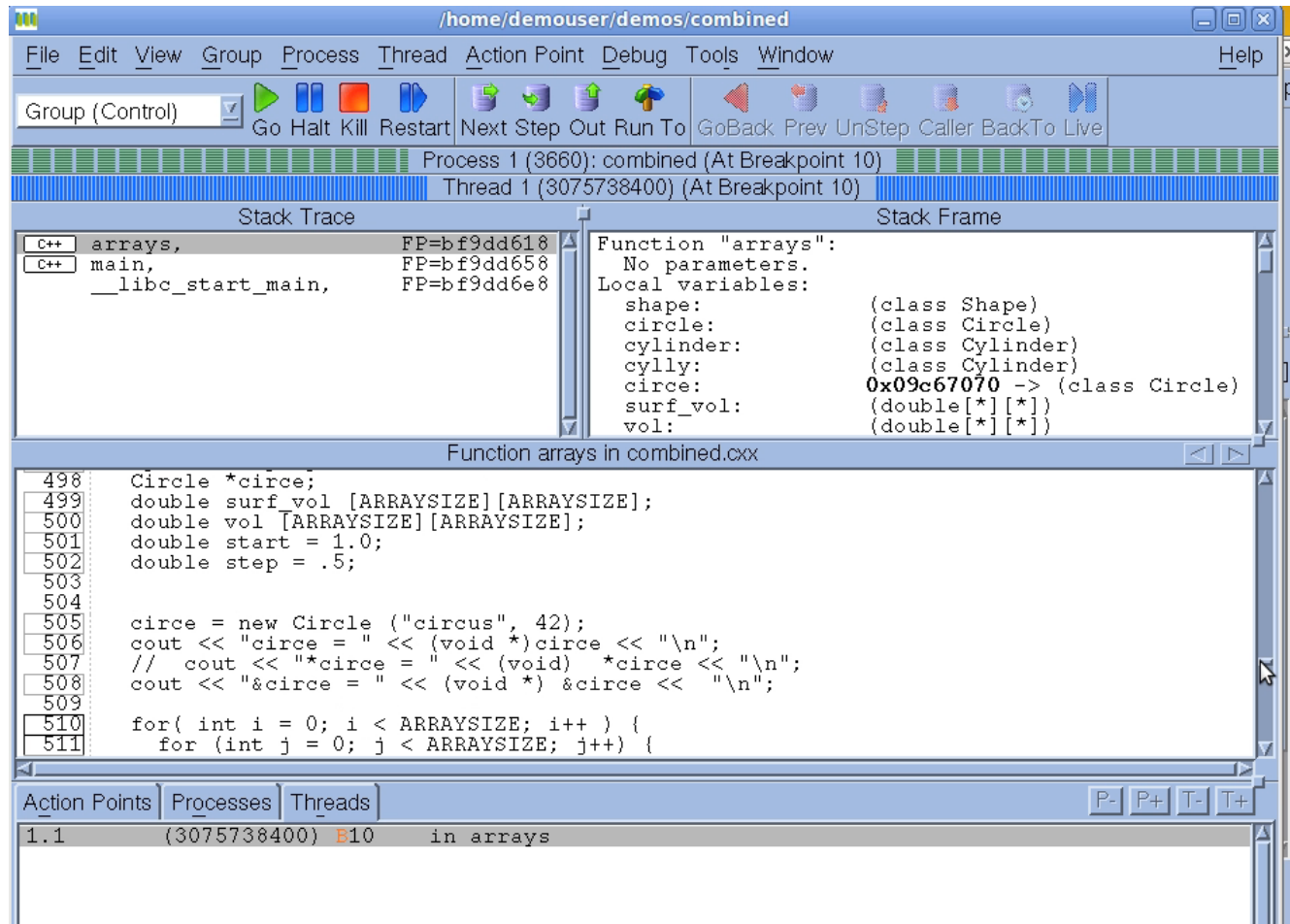  - **Find deadlocks**

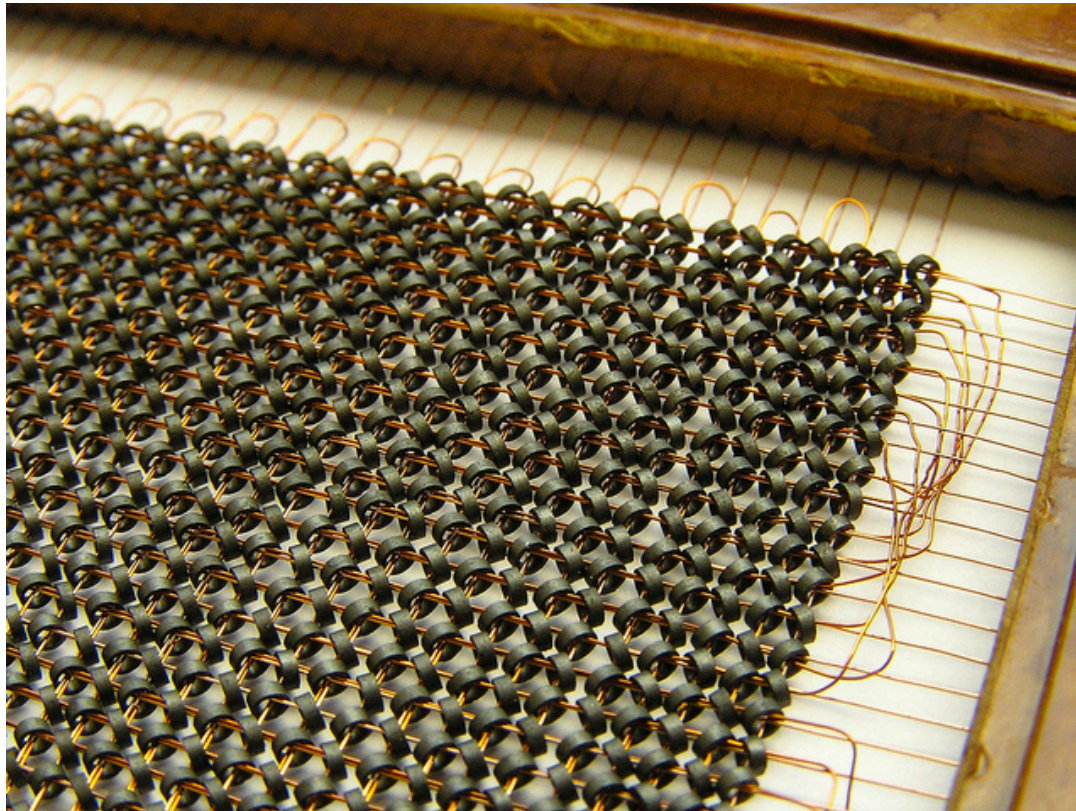# Visualize

# Visualization

ROGUE WAVE
S O F T W A R E

# Visualization

**Get the big picture – Observe anomalies – Utilize Pattern recognition – Save time!**

# … And Don't Forget the Memory!

# MemoryScape

**Memory bugs often go undetected until the worst possible time**

- Symptoms often surface long after the actual damage is done
- Some only surface after hours or even days of operation
- In many cases, the programs affected are "innocent bystanders"

**MemoryScape: Fully Integrated in TotalView**

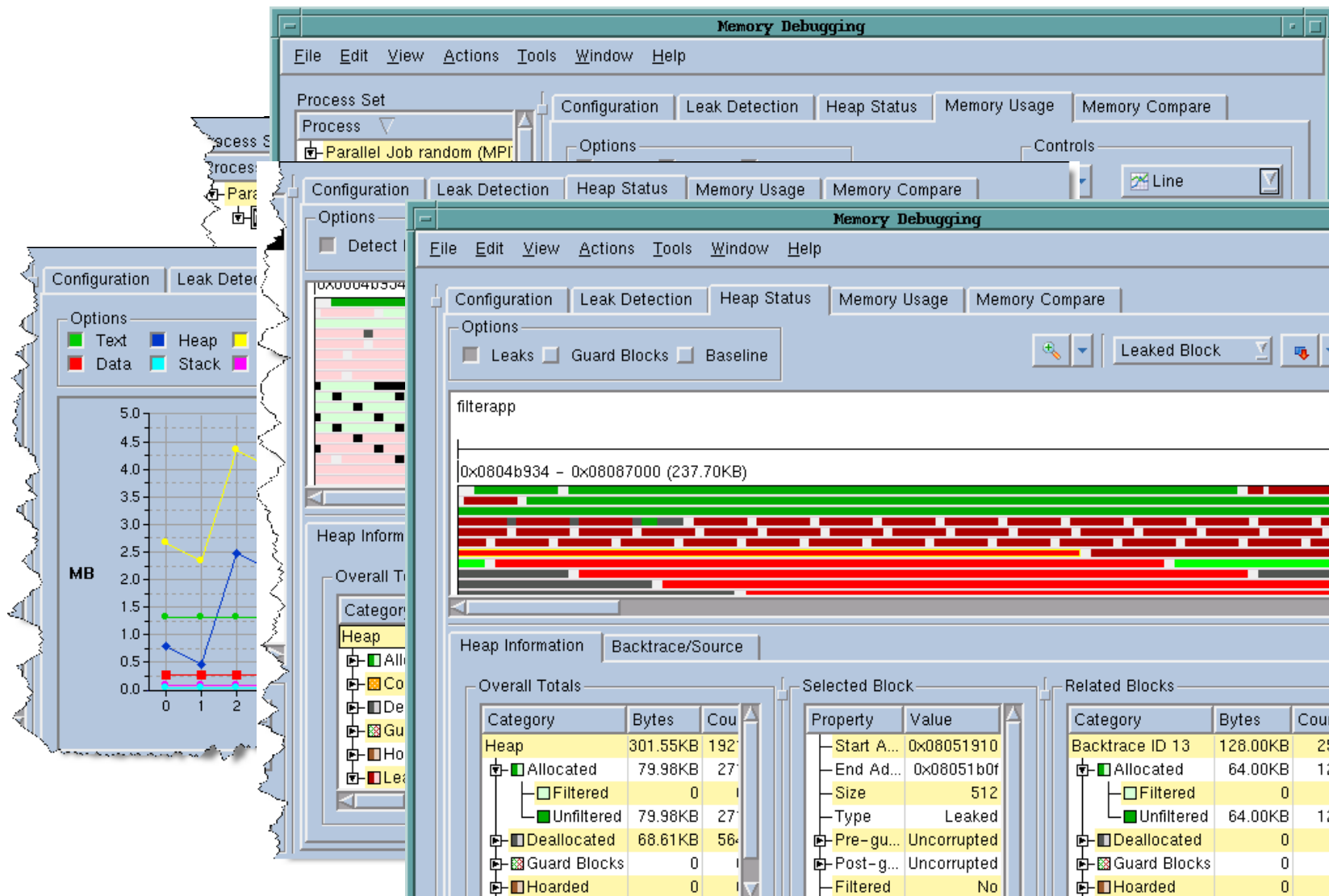**No Source Code or Binary Instrumentation**

- Use it with your existing builds
- Programs run nearly full speed
- Low performance overhead
- Low memory overhead • Efficient memory usage

**MemoryScape Feature Highlights**

- Automatic allocation problem detection
- Heap Graphical View
- Leak detection
- Block painting
- Dangling pointer detection
- Deallocation/reallocation notification
- Memory Corruption Detection - Guard Blocks
- Memory Hoarding
- Memory Comparisons between processes
- Collaboration features

# MemoryScape

# What's Coming

- ## Increased Scalability
  - **Leveraging TotalView's Architecture**
  - **Efficient Use of Cluster Resources**
    - Extremely light weight debug agents; Minimal memory footprint
    - More space on the compute nodes for user application code
  - **Tree-Based Overlay Network**
    - Broadcast of Operations; Aggregation of Events and Data

- ## Advanced User Interface
  - New GUI Framework
  - Changes focused on extreme scale debugging

- ## CUDA 4.1 now; 4.2 and 5.0 this year

- ## Replay Enhancements
  - **Record on Demand (in Beta)**
  - **Replay Debug from Core File**

- ## OpenACC Support

- ## Intel MIC Support
  - **Come see a demo at ISC '12**

# Developing for Parallel Architectures

**TotalView®**

- Code debugging
  - Highly scalable interactive GUI debugger
    - Easy to use -- without sacrificing detail that users need to debug
    - Used from workstations to the largest supercomputers
  - Powerful features for debugging multi-threaded, multi-process, and MPI parallel programs
  - Compatible with wide variety of compilers across several platforms and operating systems
- Memory Debugging
  - Parallel memory analysis and error detection
    - Intuitive for both intensive and infrequent users
  - Easily integrated into the validation process
- Reverse Debugging
  - Parallel record and deterministic replay within TotalView
    - Users can run their program "backwards" to find bugs
  - Allows straightforward resolution of otherwise stochastic bugs
- GPU CUDA Debugging
  - Full Hybrid Architecture Support
  - Asynchronous Warp Control
  - Multi-Device and MPI Support

**ROGUE WAVE®**
SOFTWARE

_Thanks!_

# ROGUE WAVE®
## S O F T W A R E

**Developing parallel, data-intensive applications is hard.**
## We make it easier.

**www.roguewave.com**